

## Assignment II – CS 6474 Social Computing\*

<i>Grade</i>	Max 75 points; 15% of overall grade (late policy applies)
<i>Due</i>	Dec 13, 2016, 11:59pm Eastern Time
<i>What to hand in</i>	A report (as a pdf file) with answers to the different questions; students choosing option A also need to include their code as a zipped folder
<i>Where to submit</i>	T-Square

### Tasks (Choose ONE of the options A or B)

**Option A.** This option of the assignment tests your design skills in addressing some of the common challenges that engender many social computing systems.

**Part 1 (Selective Exposure):** This 2016 Presidential election season invited significant media and press attention around the issue of dissemination of fake news via social media, and relatedly, that of the constraints of being inside a filter bubble<sup>1</sup>. Building on our class readings on Polarization and Selective Exposure and taking Facebook as an example social platform, answer the following two questions:

- (20 points) Include some sketches and low fidelity prototypes of a design feature (internal or external to Facebook) to allow its users bust the filter bubble, or reduce selective exposure. Describe how this feature will work, for example, what information about the end users' Facebook activity would it use, what interactions would it allow, and the extent of agency it will allow to the end users to curate their bubbles.
- (15 points) Propose a study design, of any type, through which you will assess if the above feature, if incorporated, will be successful to helping users reduce selective exposure to information. Provide your rationale behind this design, the challenges in execution of the study, and possible ways to mitigate these challenges. Assume you have unlimited time and resources to conduct this study, however you are external to Facebook, i.e., do not have access to internal policies, curation practices, or server data.

**Part 2 (Moderation and Regulation of Behaviors):** Please read the chapter titled “Regulating Behavior in Online Communities” by Sara Kiesler, Robert Kraut, Paul Resnick and Aniket Kittur [1]. In the chapter, Kiesler et al. provide an elaborate list of various design choices that could be utilized to regulate or curb non-normative or bad behavior in online communities. The table in pgs. 36-37 gives a summary of these design considerations, categorized into several groups – 1) Selection, sorting, highlighting; 2) Community Structure; 3) Feedback and Rewards; 4) Access Controls; 5) Roles, rules, policies, and procedures; and 6) Presentation and framing. Building on our readings on Deviant Communities, and Reputation, Social Signaling and Moderation, answer the following questions, focusing on three scenarios where non-normative behavior is present, and how Kiesler et al’s design choices may help regulation:

- (10 points) *Case #1: The goal is to limit pro-harassment attitudes and behaviors (e.g., around gender, sexual orientation, or race/ethnicity) on Twitter.* Discuss which of the six design choice categories will be the most appropriate here. Provide your rationale behind your choice. How does it contrast with relatively naïve regulation/moderation strategies like removals or bans? Describe how robust is your chosen design choice category – is it easy to be “gamed”; does it need considerable hands-on involvement of the moderators; could it hinder community growth and participation?
- (10 points) *Case #2: The goal is to limit flame wars in YouTube comment threads.* Discuss which of the six design choice categories will be the most appropriate here. Provide your rationale behind your choice. How does it contrast with relatively naïve regulation/moderation strategies like removals or bans? Describe how robust is

---

\* We try very hard to make questions as unambiguous as possible. If confused, send the instructor and TA an email stating the cause of confusion and your assumptions explicitly. All questions regarding confusion must be asked before 72 hours of the due date.

<sup>1</sup> <http://nymag.com/scienceofus/2016/11/how-facebook-and-the-filter-bubble-pushed-trump-to-victory.html>

your chosen design choice category – is it easy to be “gamed”; does it need considerable hands-on involvement of the moderators; could it hinder community growth and participation?

- c) (10 points) *Case #3: The goal is to limit information sharing on Reddit around behaviors that are damaging to health and well-being.* Discuss which of the six design choice categories will be the most appropriate here. Provide your rationale behind your choice. How does it contrast with relatively naïve regulation/moderation strategies like removals or bans? Describe how robust is your chosen design choice category – is it easy to be “gamed”; does it need considerable hands-on involvement of the moderators; could it hinder community growth and participation?
- d) (10 points) *Comparison of Design Choices:* Were the design choices you suggested for cases #1, #2, and #3 largely similar or dissimilar? In either case, describe reasons driving these similarities or dissimilarities. Comment on the opportunities and challenges of adopting similar or dissimilar designs in regulating different types of non-normative or bad behaviors.

**Reference:**

- [1] Kiesler, S., Kraut, R., Resnick, P., & Kittur, A. (2012). Regulating behavior in online communities. Building Successful Online Communities: Evidence-Based Social Design. MIT Press, Cambridge, MA. [\[Link to pdf\]](#)

**Option B.** The goal of this option of the assignment is to develop different supervised learning models to identify success or failure of altruistic requests on social media. The questions derive from social computing research that aims to understand linguistic markers of altruism as described on social media [1]. The questions in the assignment will test your understanding of theoretical notions of language and help seeking (narratives, moral foundations) and to what extent they can provide insights into the social construct of altruistic requests.

**Part 1:** Please refer to the enclosed zipped folder that contains dataset and associated information<sup>2</sup>. The dataset, named the file `pizza_request_dataset.json`, contains a collection of 5671 textual requests for pizza from the Reddit community “Random Acts of Pizza”<sup>3</sup> (henceforth referred to as ROAP) together with their outcome (successful/unsuccessful) and meta-data. All requests ask for the same altruistic request: a free pizza, and span the timeframe December 8, 2010 to September 29, 2013. The outcome of each request – whether its author received a pizza (successful) or not (unsuccessful) – is known. In the questions below, the ground truth data for all of the classification models will be this outcome, specifically in the file `pizza_request_dataset.json`, the field `requester_received_pizza`. Please refer to Appendix I of this assignment document for an elaborate listing and description of all of the fields in the dataset file.

The features to be used in the classification models are described in the questions below. Please develop one classifier, specifically a Support Vector Machine model with a linear kernel and default parameters corresponding to each question below. For all of the classifiers, use a randomly sampled 10% of the dataset as test set (567 posts), and the remaining 90% as the training dataset (5104 posts) – the training and test sets need to be consistent across all classifiers below, i.e., the same 567 posts should be used for testing and the same 5104 for training for a), b), c) and d).

- a) *Model 1 –  $n$ -grams (10 points):* This model will extract the top 500 unigrams and top 500 bigrams<sup>4</sup> as features to classify posts that would be successful or those that will be unsuccessful in their pizza requests. Here “top” means most frequently occurring unigrams and bigrams in the posts belonging to the training set. Using these  $n$ -gram features, train and test an SVM classifier as described above. Report a table containing the accuracy of your classifier, precision, recall, F1, specificity, and AUC.
- b) *Model 2 – Activity and Reputation (10 points):* This model will utilize a variety of the activity and reputation data included in the dataset file (`pizza_request_dataset.json`) as features to distinguish between successful and unsuccessful requests. The specific activity features will use the values included in the following fields corresponding to each post:

```
post_was_edited
requester_account_age_in_days_at_request
requester_account_age_in_days_at_retrieval
requester_days_since_first_post_on_raop_at_request
requester_days_since_first_post_on_raop_at_retrieval
requester_number_of_comments_at_request
requester_number_of_comments_at_retrieval
requester_number_of_comments_in_raop_at_request
requester_number_of_comments_in_raop_at_retrieval
requester_number_of_posts_at_request
requester_number_of_posts_at_retrieval
requester_number_of_posts_on_raop_at_request
requester_number_of_posts_on_raop_at_retrieval
requester_number_of_subreddits_at_request
requester_subreddits_at_request
```

And the specific reputation features will use the values included in the following fields for each post:

```
number_of_downvotes_of_request_at_retrieval
```

---

<sup>2</sup> Downloaded from the SNAP Stanford website: <http://snap.stanford.edu/data/web-RedditPizzaRequests.html>

<sup>3</sup> [https://www.reddit.com/r/Random\\_Acts\\_Of\\_Pizza/](https://www.reddit.com/r/Random_Acts_Of_Pizza/) Excerpt from the subreddit description: “Feel like giving a random redditor a free pizza, but don't know how or who? Well this is the right place for you! Random giving is why we are here!”

<sup>4</sup> Post content is given in the field “`request_text`” in the dataset file `pizza_request_dataset.json`.

```
number_of_upvotes_of_request_at_retrieval
requester_upvotes_minus_downvotes_at_request
requester_upvotes_minus_downvotes_at_retrieval
requester_upvotes_plus_downvotes_at_request
requester_upvotes_plus_downvotes_at_retrieval
requester_user_flair
```

Using these values for activity and reputation as features, train and test an SVM classifier as described above.

Report a table containing the accuracy of your classifier, precision, recall, F1, specificity, and AUC.

- c) *Model 3 – Narratives (15 points)*: This third model will extract features corresponding to the narrative dimensions identified in [1]. Refer to the enclosed files within “/resources/narratives”. There are five narratives – *desire, family, job, money, and student*. Each narrative file has a set of words associated with it. To extract post features corresponding to a narrative, perform regular expression match between all words corresponding to the narrative and those corresponding to a post (in the training and test sets)<sup>3</sup>. The narrative features for a post will be the ratio of the number of matches for each narrative to the total number of white spaced words in the post. Using these five narrative features, train and test an SVM classifier as described above. Report a table containing the accuracy of your classifier, precision, recall, F1, specificity, and AUC.
- d) *Model 4 – Moral foundations (15 points)*: This third model will use the dimensions of “moral foundations” as features for classifying successful and unsuccessful requests. These dimensions are based on the moral foundations theory<sup>5</sup> that seeks to understand why morality varies so much across cultures yet still shows so many similarities and recurrent themes. In brief, the theory proposes that several innate and universally available psychological systems are the foundations of “intuitive ethics.” The dimensions of the moral foundations include: *care/harm, loyalty/betrayal, authority/subversion, and sanctity/degradation*. Their descriptions can be found in Appendix II. To extract features corresponding to the different dimensions, first refer to the enclosed file “MoralFoundations.dic” under “/resources” – the file opens with any simple plain text editor program. The dictionary contains terms indexed by integers, where the integers are mapped to the moral foundations dimensions. Then, for a given post in your training or test data<sup>3</sup>, obtain one feature corresponding to each dimension, by matching (with regular expressions) each word in the dictionary for that dimension to each word in the post. This way, you will obtain a count variable of the occurrence of the dimension in the post. By dividing this count by the total number of white spaced words in the post, you will obtain a normalized feature value for the same dimension. Using these dimensions as features, train and test an SVM classifier as described above. Report a table containing the accuracy of your classifier, precision, recall, F1, specificity, and AUC.

**Part 2:** Present a discussion of the performance of the above four models:

- (2 points) Which of the four classifiers performed the best; which one performed the worst?
- (3 points) Describe your anticipated reasoning driving these differences in performance of the classifiers.
- (5 points) For models 3 and 4 in particular, describe their performance compared to models 1 and 2. Why do you think they perform better or worse than models 1 and 2? Between models 3 and 4, which one is better? What could be the reason behind this observation?
- (5 points) Present your reasoning if your models indicate that language is able to predict success of altruistic requests – other than model 2, all of the other models rely on language.

**Part 3:** Presentation a comparative discussion of the performance of all of your classification models and the performance metrics (AUC) reported in Table 4 of [1]:

- (5 points) In what ways are your models similar or different from those in Table 4 of [1]?
- (5 points) Where and why do they perform better or worse compared to [1]?

**Reference:**

- [1] Althoff, T., Danescu-Niculescu-Mizil, C., & Jurafsky, D. (2014). How to ask for a favor: A case study on the success of altruistic requests. In Proc. ICWSM 2014. [\[Link to pdf\]](#)

---

<sup>5</sup> <http://moralfoundations.org/>

## Appendix I

Format of the file pizza\_request\_dataset.json for Option A:

Field	Description
giver_username_if_known	Reddit username of giver if known, i.e. the person satisfying the request ("N/A" otherwise).
in_test_set	Boolean indicating whether this request was part of our test set.
number_of_downvotes_of_request_at_retrieval	Number of downvotes at the time the request was collected.
number_of_upvotes_of_request_at_retrieval	Number of upvotes at the time the request was collected.
post_was_edited	Boolean indicating whether this post was edited (from Reddit).
request_id	Identifier of the post on Reddit, e.g. "t3_w5491".
request_number_of_comments_at_retrieval	Number of comments for the request at time of retrieval.
request_text	Full text of the request.
request_text_edit_aware	Edit aware version of "request_text". We use a set of rules to strip edited comments indicating the success of the request such as "EDIT: Thanks /u/foo, the pizza was delicious".
request_title	Title of the request.
requester_account_age_in_days_at_request	Account age of requester in days at time of request.
requester_account_age_in_days_at_retrieval	Account age of requester in days at time of retrieval.
requester_days_since_first_post_on_raop_at_request	Number of days between requesters first post on RAOP and this request (zero if requester has never posted before on RAOP).
requester_days_since_first_post_on_raop_at_retrieval	Number of days between requesters first post on RAOP and time of retrieval.
requester_number_of_comments_at_request	Total number of comments on Reddit by requester at time of request.
requester_number_of_comments_at_retrieval	Total number of comments on Reddit by requester at time of retrieval.
requester_number_of_comments_in_raop_at_request	Total number of comments in RAOP by requester at time of request.
requester_number_of_comments_in_raop_at_retrieval	Total number of comments in RAOP by requester at time of retrieval.
requester_number_of_posts_at_request	Total number of posts on Reddit by requester at time of request.
requester_number_of_posts_at_retrieval	Total number of posts on Reddit by requester at time of retrieval.
requester_number_of_posts_on_raop_at_request	Total number of posts in RAOP by requester at time of request.
requester_number_of_posts_on_raop_at_retrieval	Total number of posts in RAOP by requester at time of retrieval.
requester_number_of_subreddits_at_request	The number of subreddits in which the author

requester_received_pizza	had already posted in at the time of request.
requester_subreddits_at_request	Boolean indicating the success of the request, i.e., whether the requester received pizza.
requester_upvotes_minus_downvotes_at_request	The list of subreddits in which the author had already posted in at the time of request.
requester_upvotes_minus_downvotes_at_retrieval	Difference of total upvotes and total downvotes of requester at time of request.
requester_upvotes_plus_downvotes_at_request	Difference of total upvotes and total downvotes of requester at time of retrieval.
requester_upvotes_plus_downvotes_at_retrieval	Sum of total upvotes and total downvotes of requester at time of request.
requester_user_flair	Sum of total upvotes and total downvotes of requester at time of retrieval.
requester_username	Users on RAOP receive badges (Reddit calls them flairs) which is a small picture next to their username. In our data set the user flair is either None (neither given nor received pizza, N=4282), "shroom" (received pizza, but not given, N=1306), or "PIF" (given after received, N=83).
unix_timestamp_of_request	Reddit username of requester.
unix_timestamp_of_request_utc	Unix timestamp of request (supposedly in timezone of user but in most cases equal to the UTC timestamp which is incorrect since most RAOP users are from the USA).

## **Appendix II**

Descriptions of the different moral foundations dimensions:

*Care/harm:* This foundation is related to our long evolution as mammals with attachment systems and an ability to feel (and dislike) the pain of others. It underlies virtues of kindness, gentleness, and nurturance.

*Fairness/cheating:* This foundation is related to the evolutionary process of reciprocal altruism. It generates ideas of justice, rights, and autonomy.

*Loyalty/betrayal:* This foundation is related to our long history as tribal creatures able to form shifting coalitions. It underlies virtues of patriotism and self-sacrifice for the group. It is active anytime people feel that it's "one for all, and all for one."

*Authority/subversion:* This foundation was shaped by our long primate history of hierarchical social interactions. It underlies virtues of leadership and followership, including deference to legitimate authority and respect for traditions.

*Sanctity/degradation:* This foundation was shaped by the psychology of disgust and contamination. It underlies religious notions of striving to live in an elevated, less carnal, more noble way.